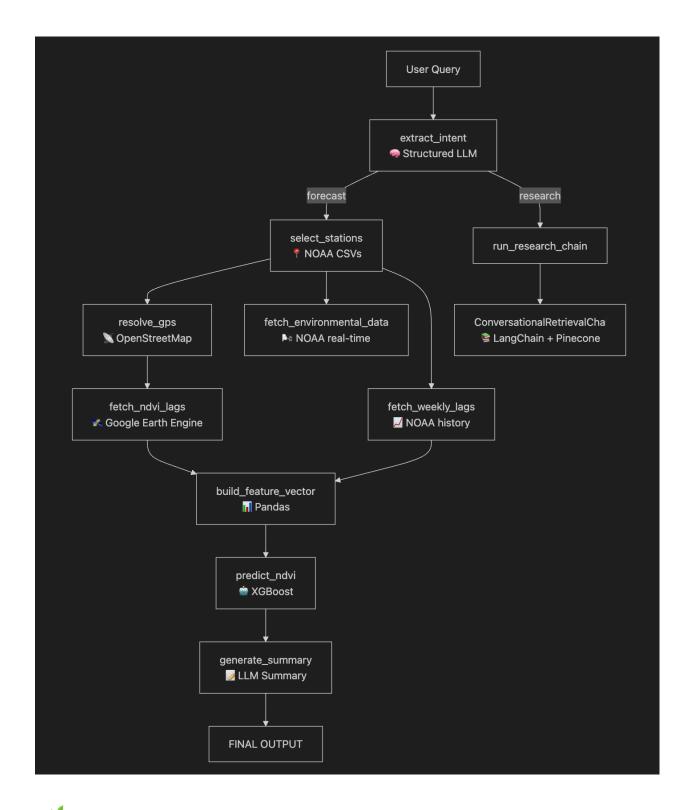
## **LangGraph Pipeline**

<ul><li>Oreated</li></ul>	@March 27, 2025 6:20 PM
<u>≔</u> Tags	



# Mangrove Intelligence Agent – System Manual

### **Overview**

This system is a hybrid AI agent built with **LangGraph** and **LangChain**. It intelligently handles **two types of user questions**:

- 1. Forecast questions real-time, location-specific mangrove health analysis
- Research questions general, factual, or conceptual questions about mangroves

### Agent Behavior Logic

When the user asks a question:

- 1. The agent first uses an LLM to **extract intent**:
  - Is it a forecast or research query?
  - Does the question mention a location?
  - Can we infer a U.S. state from that location?
- 2. Then, depending on the detected **goal**, it either:
  - Activates a LangGraph pipeline (for forecast)
  - Routes to a LangChain retrieval-based agent (for research)

### Forecast Agent (LangGraph)

#### ➤ Goal

Deliver a 5–6 sentence summary of current and predicted mangrove conditions using:

- Real-time data (wind & water)
- Satellite imagery (NDVI from Earth Engine)
- Machine learning prediction (XGBoost)

### 

#### 1. extract\_intent\_node

- Parses user\_query
- Returns: goal , location , state
- **LangGraph Workflow: Routing** (decides forecast VS research path)

### 2. Router (Inside LangGraph)

- If goal == forecast → continue
- If goal == research → exit graph with flag
- ZangGraph Workflow: Routing

#### 3. select\_stations

- Loads NOAA station metadata from CSVs (per state)
- Uses LLM to rank best wind/water stations near given location

### 4. Parallel Data Fetching

These run in parallel after station selection:

- resolve\_gps: Gets latitude/longitude from location
- fetch\_environmental\_data: Pulls real-time wind/water values from NOAA API
- fetch\_weekly\_lags: Gets 8 weeks of historical hourly data from NOAA, converts to weekly averages

#### 5. fetch\_ndvi\_lags

- Uses Google Earth Engine to extract 8 weeks of NDVI data from MODIS imagery

#### 6. build\_feature\_vector

- Combines all time series into a unified DataFrame
- Adds lag features for wind, tide, and NDVI

### 7. predict\_ndvi

• Applies scaler + XGBoost model to forecast NDVI for current week

### 8. generate\_summary

- Uses LLM to summarize:
  - Wind speed
  - Water level
  - NDVI prediction
  - Original user intent
- **LangGraph Workflow: Prompt Chaining** (LLM + context summary)

### **External Services Used**

Node	Source / API
resolve_gps	OpenStreetMap Nominatim
fetch_environmental_data	NOAA Tides & Currents API
fetch_weekly_lags	NOAA Historical API (chunked)
fetch_ndvi_lags	Google Earth Engine (MODIS)
build_feature_vector	Pandas (time series wrangling)
predict_ndvi	XGBoost, Scikit-learn Scaler
generate_summary	LangChain Prompt Template + LLM

### Research Agent (LangChain)

### ➤ Goal

Answer open-ended, non-location-specific questions about mangroves

#### **Uses:**

- LangChain's ConversationalRetrievalChain
- A Pinecone vector index of documents

- OpenAl Embeddings
- Memory to keep track of chat history

### **Example queries:**

- "Why are mangroves important?"
- "What are common threats to mangroves?"

### Modular Roles

Component	Purpose
extract_intent_node	Centralized logic for goal inference
run_agent()	Top-level function that wraps LangGraph and LangChain agents
forecast_chain()	Invokes LangGraph pipeline for forecasts
run_research_chain()	Invokes LangChain retriever agent
memory	Stores user & assistant messages across turns

### Example Usage

response = run\_agent("How are mangroves doing in Key West?")
# → Returns a real-time forecast summary using LangGraph

response = run\_agent("What are mangroves?")
# → Returns an educational explanation using Pinecone + RAG

### P Developer Notes

- All forecast logic is encapsulated in LangGraph and benefits from:

  - Prompt Chaining
- Easy debugging via StateGraph.get\_graph().draw\_mermaid()

• LangChain is used **only** for conversational memory + document retrieval

### Extensibility Ideas

Idea	How
Add human-in-the-loop review	Use LangGraph's pause / memory checkpoint
Add streaming responses	Use LangGraph streaming + token events
Add more data sources	Plug in extra nodes (salinity, temperature, etc.)
Allow follow-up clarification	Use LangChain memory thread

### 7. continuous\_evaluator (Optional)

- Inputs: NDVI prediction + Actual NDVI (from delayed source)
- Loop:
  - Evaluate accuracy
  - If bad, log feedback → optional retraining
- LangGraph Pattern: Evaluator-Optimizer

```
flowchart TD

%% —— ENTRY POINT ——

A[User Query] → B[extract_intent<br/>
%% —— ROUTING ——

B → |forecast| C[select_stations<br/>
B → |research| Z[run_research_chain]

%% —— PARALLEL FETCH AFTER STATION SELECTION ——

C → C1[resolve_gps<br/>
by OpenStreetMap]

C → C2[fetch_environmental_data<br/>
c → C3[fetch_weekly_lags<br/>
c → C3[fetch_weekly_lags<br/>
c → NDVI PARALLEL BRANCH ——

C1 → D1[fetch_ndvi_lags<br/>
br/>
Google Earth Engine]

%% —— FAN-IN FOR FEATURE VECTOR ——

C3 → E[build_feature_vector<br/>
li Pandas]

D1 → E
```

```
%% ——— PREDICTION AND SUMMARY ———
E → F[predict_ndvi<br/>
\Rightarrow XGBoost]
F → G[generate_summary<br/>
\Rightarrow LLM Summary]
G → H[FINAL OUTPUT]
%% ——— RESEARCH SIDE PATH ———
Z → ZZ[ConversationalRetrievalChain<br/>
\Rightarrow LangChain + Pinecone]
```

### **Details**

### **Model Flow and Implementation Summary**

### Extract goal and location from user query

- extract\_intent\_node
- Parses goal, location, state
- Stores into state for downstream use

### Goal Branch: summary or qa

- ✓ Supports "summary" → full forecast pipeline
- X "qa" not yet implemented (planned in diagram)

### Find nearby stations

- select\_station\_by\_location\_node
- Uses LLM to rank NOAA stations by proximity
- Outputs station\_ids and station\_candidates

#### **Resolve GPS coordinates**

- resolve\_gps\_from\_location\_node
- Geocodes location, state → outputs gps tuple

#### Fetch real-time wind and water data

- fetch\_environmental\_data\_node
- Pulls latest values from NOAA API
- Output: environmental\_data dict

#### Fetch 8-week wind and water history

- Performed inside build\_feature\_vector\_node
- Uses fetch\_weekly\_noaa\_lags\_chunked(...)

#### Fetch 7-week NDVI history from MODIS

- Also inside build\_feature\_vector\_node
- Uses get\_cleaned\_weekly\_ndvi\_series(...)

#### **Build feature vector from time series**

- build\_feature\_vector\_node
- Constructs 23-dimensional input from current + lagged values
- Output: feature\_vector + feature\_df (for diagnostics)

### **Predict NDVI using XGBoost model**

- predict\_ndvi\_node
- Scales inputs using scaler, runs prediction
- Output: ndvi\_prediction (float)

### **Generate natural language summary**

- generate\_summary\_node
- Interprets NDVI result and returns textual summary

#### Return final result

- V Final State includes:
  - User input

- Station metadata
- Real-time & historical data
- Feature vector
- Prediction
- Summary text

### **QA Goal Handling**

- Ont yet implemented
- Can be added as a simple LLM QA branch